



Home



Search



List

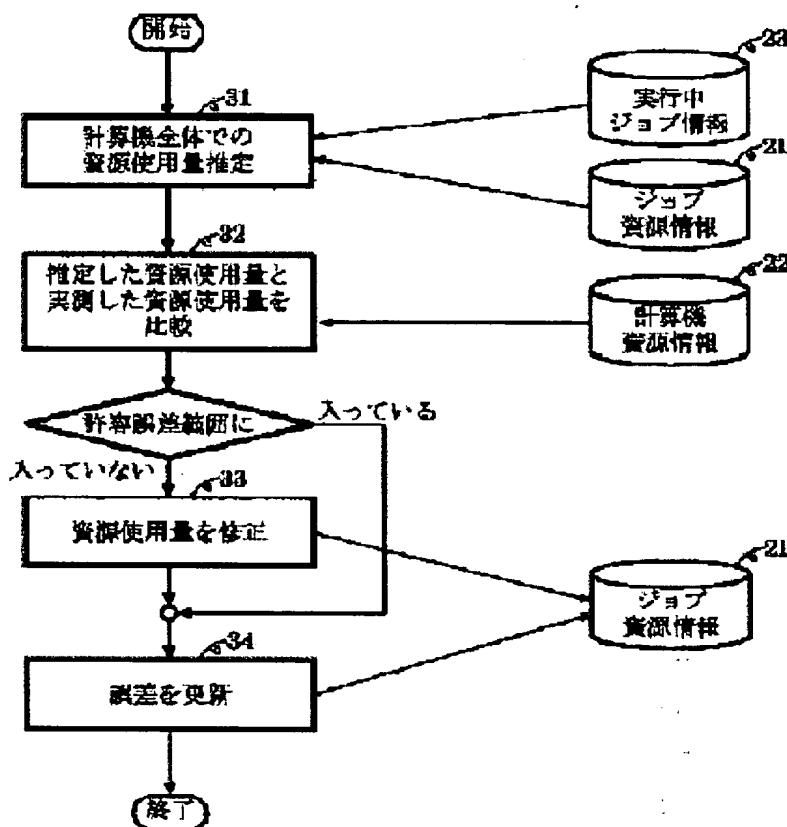
☐ Include

MicroPatent® PatSearch FullText: Record 1 of 1

Search scope: US Granted JP (bibliographic data only)

Years: 1991-2003

Patent/Publication No.: JP09311795



Order This Patent

Family Lookup

Find Similar

Legal Status

[Go to first matching text](#)

JP09311795 A
SCHEDULING METHOD
HITACHI LTD

Inventor(s): KURODA SAWAKI; FUKUZAWA JUNJI; YAMAMURA NOBUO; MORITA SHINJI
Application No. 08126847 JP08126847 JP, Filed 19960522, A1 Published 19971202

Abstract: PROBLEM TO BE SOLVED: To automatically acquire a system resource quantity needed to execute individual jobs by updating the difference between a measured value of the resource consumption of a computer and an estimated value and error in the resource demand quantity of a job being executed.

SOLUTION: At each resource, job information 32 in execution and job resource information 21 are totalized to find the estimated value of the resource consumption of the whole computer (step 31), and the found estimated value of the resource consumption is compared with the measured value of the consumption of the resource information 22

of the computer (step 32). When an error in the resource consumption of the whole computer is exceeded, the resource consumption is corrected according to the rate of the resource demand quantity and error (step 33), and the job resource information in execution and further the difference between the measured value of the resource consumption of the computer and the estimated value of the resource consumption and the error in the demand quantity of the job resource information 21 in execution are updated (step 34). Therefore, the system resource quantity needed to execute individual jobs can automatically be acquired.

Int'l Class: G06F00946; G06F01516

Patents Citing this One: No US, EP, or WO patents/search reports have cited this patent.

[Home](#)[Search](#)[List](#)

For further information, please contact:
[Technical Support](#) | [Billing](#) | [Sales](#) | [General Information](#)

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平9-311795

(43) 公開日 平成9年(1997)12月2日

| (51) Int.Cl. ⁶ | 識別記号 | 庁内整理番号 | F I | 技術表示箇所 |
|---------------------------|-------|--------|--------------|---------|
| G 0 6 F 9/46 | 3 4 0 | | G 0 6 F 9/46 | 3 4 0 F |
| | | | | 3 4 0 D |
| 15/16 | | | 15/16 | 4 2 0 J |

審査請求 未請求 請求項の数 6 O L (全 9 頁)

(21) 出願番号 特願平8-126847

(22) 出願日 平成8年(1996)5月22日

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(72) 発明者 黒田 沢希

神奈川県川崎市麻生区王禅寺1099番地株式
会社日立製作所システム開発研究所内

(72) 発明者 福澤 淳二

神奈川県川崎市麻生区王禅寺1099番地株式
会社日立製作所システム開発研究所内

(72) 発明者 山村 宣夫

神奈川県横浜市戸塚区戸塚町5030番地株式
会社日立製作所ソフトウェア開発本部内

(74) 代理人 弁理士 小川 勝男

最終頁に続く

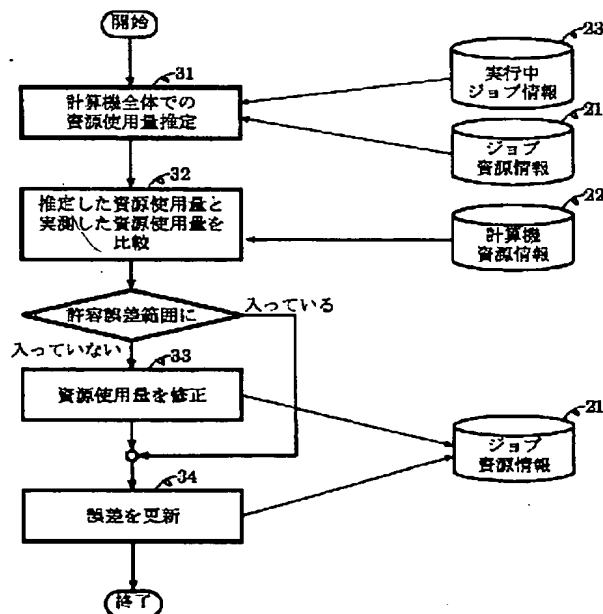
(54) 【発明の名称】 スケジュール方法

(57) 【要約】

【課題】 ジョブ投入後のシステムの状態を予測するために、ジョブの特性を自動的に収集することを目的とする。

【解決手段】 分散計算機システムにおいて、個別のジョブの使用する資源量の情報とその確度を備え、実行中のジョブの前記資源量の合計と、計算機システムの負荷の実測値との差を、使用資源量情報確度に応じてジョブの使用する資源量の情報を修正することによって、ジョブの使用する資源量の情報を獲得し、ジョブの使用する資源量を元に、特定の資源に集中しないスケジューリングを行なう。

図3



【特許請求の範囲】

【請求項 1】 計算機システムが提供する資源と、計算機システム全体での資源の使用量を測定する手段と、該計算機上で実行中のジョブと、実行中ジョブの一覧情報と、を備える計算機システムにおける、ジョブ資源要求量獲得方法であって、

ジョブの名称と、計算機の各資源について当該ジョブが要求する量の推測値と、その確からしさをあらわす確度とを含むジョブ資源情報を備え、

(a) 計算機で実行中のジョブについて、ジョブ資源情報の当該ジョブが要求する資源量を合計して計算機全体の資源使用量の推測値をもとめ、当該ジョブに関する資源要求量の誤差を加算して、計算機全体の資源使用量の誤差をもとめ、

(b) 計算機の資源使用量の実測値と、前記求めた資源使用量の推測値との差が、前記計算機全体の資源使用量の誤差より大きい時、

(c) 資源使用量の実測値と推測値の差がなくなるように、当該計算機で実行中のジョブについて、資源要求量およびその誤差の割合に応じて、実行中のジョブ管理情報の推定資源要求量に加え、

(d) 前記ステップ(b)で求めた実測値と推測値の差と、実行中のジョブの資源要求量の誤差を更新することを特徴とした、スケジュール方法。

【請求項 2】 請求項 1 の計算機システムで、さらに計算機の持つ資源容量の情報を備え、前記ステップ(a)において、

(a1) 実行中のジョブおよび各資源について、当該ジョブの資源推定要求量の合計を出力し、

(a2) 資源推定要求量合計と、計算機の資源容量と、各ジョブの推定資源要求量とからジョブの稼働率を算出し、

(a3) 各ジョブの稼働率から、ジョブの資源推定使用量と、処理時間の推定値を出力することを特徴とした、スケジュール方法。

【請求項 3】 請求項 1 の計算機システムにおいて、スケジューラがジョブの起動／終了を検知するステップを備え、前記検知した時に、前記各ステップを実行することを特徴とする、スケジュール方法。

【請求項 4】 分散計算機環境で、ジョブを投入する時、

(e) 請求項 1 で求めたジョブの資源使用量を元に、当該ジョブの資源推定使用量情報を求め、

(f) 当該ジョブの資源推定使用量を、各計算機の資源使用量に加えて、当該ジョブ投入後の資源使用量を求め、
(g) 資源使用量と資源容量の比を、各資源について和を求め、

(h) 前記和がもっとも小さくなる計算機に当該ジョブを割り当てることを特徴とする、スケジュール方法。

【請求項 5】 請求項 1 のステップ(f) の代わりに、

(f') 実行中のジョブジョブと、投入するジョブについて、資源推定使用量を読み出し、前記実行中および投入ジョブの資源推定使用量の総和を求め、これを資源使用量とすることを特徴とした、スケジュール方法。

【請求項 6】 請求項 1 のステップ(f) の代わりに、

(f'') 当該ジョブの資源推定使用量と誤差範囲の割合との積を、各計算機の資源使用量に加えて、当該ジョブ投入後の資源使用量を求める、ことを特徴とした、スケジュール方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は繰り返し実行されるジョブに関して、実行時のシステムに与えるCPU負荷やIO 負荷等から当該ジョブが必要とする負荷を推定し、当該推定値を使って、ジョブ投入後のシステムの状態を予測して、ジョブを投入する計算機を決定するスケジューラに関する。

【0002】

【従来の技術】 “PVM Guide” (http://csepl.phy.ornl.gov/pvm_guide/pvm_guide.html)によると、分散計算機環境で性能向上を図る手段として、負荷分散がある。各計算機の実行待ちジョブ数を計測し、この実行待ちジョブ数の最も少ない計算機で新しいジョブを実行することによって、各計算機の負荷を平均化することができる。

【0003】 また、Greorge A. Champine 著の『Distributed Computer Systems: Impact on Management, Design, and Analysis』によると、個別ジョブがシステムに与える負荷を得る方法として、ジョブの実行ステップ数等をエミュレータ等を利用して計測する方法がある。この方法では、ジョブの動作が詳細にわかるため、資源に与える負荷の質／量／時間をほぼ正確に取得することができる。しかし、1 命令をエミュレータで解析／実行するので、ジョブの実行ステップ数の計測には非常に時間がかかる。またジョブが他社の開発した物の場合、実行ステップ数の計測は、使用契約上認められないことがある。

【0004】 個別ジョブがシステムに与える負荷を得る別の方法として、計算機上でそのジョブだけを動かして、計算機の負荷を実測する方法がある。この方法では、短時間で負荷値を測定できる。しかし、目的のジョブ以外の影響を排除するために、測定専用計算機システムを用意する必要がある。

【0005】

【発明が解決しようとする課題】 前記システムの負荷状態による負荷分散スケジュールでは、長時間かかるジョブを同一計算機に割り当てる場合があり、その結果全体の処理時間が伸びる場合があるという問題がある。

【0006】 個別ジョブの実行負荷を得るために、当該ジョブの実行ステップ数等を計測する方法では、ジョブの実行ステップ数の計測は通常シミュレータを利用する

ため、非常に時間がかかるという問題がある。またジョブが他社の開発した物の場合、実行ステップ数の計測は、その使用契約上認められないことがあるという問題もある。

【0007】個別ジョブの実行負荷を得るために、計算機上で当該ジョブだけを動かして負荷を実測する方法では、測定専用に分散システムを用意する必要があり、コストがかかってしまうという問題がある。

【0008】本発明の目的は、個別ジョブに関して、当該ジョブの実行に必要なシステム資源量を自動的に獲得することと、前記資源量を元にジョブを計算機に割り当てることである。

【0009】

【課題を解決するための手段】ジョブ名称に対応して、当該ジョブが使用するCPUやIOの負荷およびジョブ所要時間などの推定値と、その確度を備えた情報を備える。

【0010】分散した計算機毎にCPU/IO能力などの値を測定する。

【0011】未知のジョブの場合、負荷の推定値および確度は、あらかじめ定めてあるデフォルト値を設定する。

【0012】ジョブ実行後の計算機の負荷変動を、そのジョブの推定負荷として登録する。

【0013】ジョブの実行中に実行中のジョブの負荷推定値の合計と、実測した計算機の負荷との差が、確度から求められる許容誤差の範囲内であれば、確度を上げる。

【0014】ジョブの実行中に、実行中のジョブの負荷推定値の合計と、実測した計算機の負荷との差が、確度から求められる許容誤差範囲を越えた時、推定の誤差の割合に対応して、推定の負荷の大きさに比例して、その差分を分配して負荷推定値を修正し、確度を下げる。

【0015】これらのステップにより、個別ジョブが計算機システムに与える負荷を自動的に取得する。

【0016】

【発明の実施の形態】以下に本発明の一実施の形態を説明する。

【0017】図1は本発明によるスケジュール方式を採用した分散計算機システム(1)の概要を示す。この分散計算機システム(1)は、管理計算機(11)と、複数の計算サーバ(12)と、それら繋ぐネットワーク(13)とを含む。なお、管理計算機(11)は計算サーバ(12)の一つを兼用しても構わない。

【0018】管理計算機(11)は、ジョブ資源情報(21)と、当該計算機システムに要求された起動待ちのジョブ(112)を管理するジョブ受付キュー(111)と、当該ジョブを1つの計算サーバ(12)に割り当てるスケジューラ(113)を備える。

【0019】計算サーバ(12)は、CPU時間、メモリ容

量、ファイルアクセス頻度、ネットワークアクセス頻度など、複数の資源(121)と、資源の容量や使用量を保持する計算機資源情報(22)と、資源を使用するジョブ(123)と、ジョブを管理する実行中ジョブ情報(23)と、計算サーバ(12)の資源(121)の状態を調べる資源使用量測定手段(122)と、ジョブ情報獲得手段(124)を備える。

【0020】ジョブ情報獲得手段(124)は、資源使用量測定手段(122)から計算サーバ(12)の資源の使用量と、実行中ジョブ情報(23)から、実行中のジョブ(123)が使用する資源(121)を推測し、ジョブが使用する資源の量を記憶するジョブ資源情報(12)を更新する。

【0021】ジョブの実行要求があると、実行要求受付キュー(111)に実行待ちジョブ(112)を入れる。スケジューラ(113)は実行要求受付キュー(111)から実行待ちジョブ(112)を取りだし、ジョブ資源情報(21)からジョブが使用する資源量の情報を取得し、各計算サーバ(12)の資源使用量測定手段(122)から資源の状態を取得し、実行待ちジョブ(112)を適当な計算サーバ(12)に割り当てる。

【0022】図2は、各テーブルの構造を示している。

【0023】ジョブ資源情報(21)は、分散計算機システム上に1つあり、ジョブ名(211)と、当該ジョブの実行時に消費する各資源名(212)と、その資源を当該ジョブが使用する資源必要量(213)とその確からしさをあらわす確度(214)とを備える。資源としてジョブの起動から終了までの処理時間、単位時間当たりを使用するCPU時間、占有するメモリ量、ファイルアクセス頻度、ネットワークアクセス頻度などがある。このテーブルはジョブ情報獲得手段(124)が後述の方法によって追加/変更する。

【0024】計算機資源情報(22)は、各計算サーバにあり、当該計算機が提供できる各資源について、その資源名(221)と、提供できる能力をあらわす資源容量(222)と、資源の使用が著しく悪くなる使用率をあらわす使用限界(223)と、資源に対する当該計算サーバにおける負荷の測定値である資源使用量(224)を備える。このテーブルの資源使用量(224)以外の項目は、別の手段によってあらかじめ値が入っているものとする。資源使用量(224)は、資源使用量測定手段(122)によって定期的に更新する。

【0025】実行中ジョブ情報(23)は、実行中のジョブ(123)について、ジョブ名(231)と、ジョブID(232)を備える。新規起動したジョブについて、当該ジョブのためのエントリを追加し、終了したジョブに関しては、当該終了したジョブのエントリを削除する。

【0026】以下では、前記の構成を持つ分散計算機システムにおけるジョブ情報獲得方法とスケジュール方法について、説明する。

【0027】図3はジョブ情報獲得手段におけるジョブ

情報獲得方法の詳細を示す。

【0028】定期的に以下の処理をおこなう。

【0029】ステップ31では、各資源について、実行中ジョブ情報(23)のそれぞれのジョブ名(231)と、ジョブ資源情報(21)のジョブ名(211)と一致するエントリから、から当該ジョブの資源要求量(213)を読みだす。

【0030】実行中のジョブ計算機上で複数のジョブが

$$LA_{ms} = \min \left(LC_{ms}, \sum_p LR_{ps} \right) \dots\dots\dots (数1)$$

【0033】とし、計算機mでのジョブpの稼働率Ampは、資源sの容量LCms、ジョブpの資源要求量LRpsを使って、

$$A_{mp} = \begin{pmatrix} \frac{LR_{ps}}{\sum LR_{ps}} & | LC_{ms} < \sum LR_{ps} \text{の時} \\ 1 & | LC_{ms} > \sum LR_{ps} \text{の時} \end{pmatrix} \dots\dots\dots (数2)$$

【0035】によって求める。

【0036】次にステップ32では、計算機資源情報(22)の実測した資源使用量(224)と、前記使用量の推定値LAmsとの差を求める。後の説明のためにこの差を実推差と呼ぶ。また、資源要求量(213)と確度(214)から求める許容誤差を求める。

【0037】許容誤差Epsは、ジョブpの必要資源量LRpsと、その確度Rpsと、を用いて、

【0038】

【数3】

$$E_{ps} = (1 - R_{ps}) \cdot LR_{ps} \dots\dots\dots (数3)$$

$$LR_{ps} \Leftarrow LR_{ps} + c \cdot A_{mp} \cdot \frac{LR_{ps} \cdot (1 - R_{ps})}{\sum_p (LR_{ps} \cdot (1 - R_{ps}))} \cdot (LS_{ms} - LA_{ms}) \dots\dots\dots (数4)$$

【0043】で更新する。

【0044】ステップ34では、実推差と、ジョブの資源必要量(213)とその確度(214)から、ジョブの資源必要量の確度を変更する。

【0045】ジョブpの資源sに関する資源必要量の確

$$R_{ps} \Leftarrow R_{ps} + c \cdot (E_{ps} - |LS_{ms} - LA_{ms}|) \cdot A_{mp} \cdot \frac{|LS_{ms} - LA_{ms}|}{LA_{ps}} \dots\dots\dots (数5)$$

【0047】で更新する。

【0048】以上の処理を各計算サーバ(12)の各資源について繰り返す。

【0049】次に図4はスケジューラのジョブ起動割当手段の詳細を示す。

【0050】ジョブ実行要求キュー(111)に実行待ちのジョブがある時、以下の処理をおこなう。

動いている時、それぞれの資源に対してジョブは互いに競合しあうため、ジョブの稼働率が下がる。

【0031】計算機mの資源sの使用量推定値LAmsは、実行中のジョブpとその資源要求量推定値LRpsと計算機mの資源容量LCmsを使って、

【0032】

【数1】

【0034】

【数2】

【0039】とする。

【0040】前記許容誤差より実推差が大きい時、ステップ33では、ジョブの実測値(224)と推定使用量LAmsとが一致するように、実推差を実行中のジョブの資源必要量に分配する。このとき、ジョブの資源要求量が大きいものがよりたくさん、確度の低いものがよりたくさん、変化するようにする。

【0041】ジョブpの必要資源量LRpsを、その確度Rpsと、計算機mでの測定値LSms、推定値LAms、定数cを用いて、

【0042】

【数4】

度Rpsを、計算機mでの資源使用量測定値LSmsと、資源使用量推定値LAmsと、許容誤差Epsと、定数cを用いて、

【0046】

【数5】

【0051】ステップ41では、実行予定ジョブをジョブ実行要求キュー(111)から取得し、ジョブ資源情報(21)からジョブ名(211)が合致するエントリの各資源について、資源必要量(213)とその確度(214)を取得する。

【0052】ステップ42では、計算機資源情報(22)から、実測した資源使用量(224)を取得し、ステップ41

で得たジョブの資源必要量(213)とその確度(214)を利用して、ジョブ投入後の資源使用量を推定する。

【0053】ステップ43では、各計算機の資源容量(22)と、前記ジョブ投入後の資源使用量推定値とから、推定される資源使用率を求め、これを使用限界(223)と比較する。

【0054】推定される資源使用率が使用限界(223)以下となる計算サーバ(12)がある時、ステップ44では、

$$0 > \min_m \left(\max_s \left(\frac{LS_{ms} + LR_{ps}}{LC_{ms}} - B_{ms} \right) \right) \dots\dots\dots (数6)$$

【0057】の条件を満たす。

【0058】前記条件を満たす計算サーバ(12)がない時、ステップ45では、当該ジョブを起動しないで、実行待ちキュー(111)に保留し、一定時間後に再びステップ41から処理する。

【0059】別のスケジューリング方法として、前記ステップ42において、実行中ジョブ情報(23)から、実行中のジョブの一覧を取得し、そのそれぞれの資源必要量を求め、ステップ41で得たジョブの資源必要量(21

$$0 > \min_m \left(\max_s \left(\frac{\sum_{p=\text{実行中}} (LR_{ps}) + LR_{ps}}{LC_{ms}} - B_{ms} \right) \right) \dots\dots\dots (数7)$$

【0062】の条件を満たす。

【0063】別のスケジューリング方法として、前記ステップ44において、全てのマシンの中で、ジョブ投入後の資源使用率の合計が最も低くなる計算機に割り当てる。

【0064】すなわち、この計算機 m は、ジョブ p 投入

$$0 > \min_m \left(\sum_s \left(\frac{LS_{ms} + LR_{ps}}{LC_{ms}} - B_{ms} \right) \right) \dots\dots\dots (数8)$$

【0066】の条件を満たす。

【0067】また別のスケジューリング方法として、前記ステップ44において、全てのマシンの中で、ジョブ投入後の最大の資源使用率の最も低くなる計算機に割り当てる。

【0068】すなわち、この計算機 m は、ジョブ p 投入時に、資源 s について、投入前の実測値 LS_{ms} と、投入するジョブの資源要求量 LR_{ps} と、資源容量 LC_{ms} と、資源使用限界 B_{ms} と、において、

【0069】

【数9】

$$\min_m \left(\max_s \left(\frac{LS_{ms} + LR_{ps}}{LC_{ms}} \right) \right) \dots\dots\dots (数9)$$

【0070】の条件を満たす。

【0071】ジョブが必要とする資源量が不確かな間は、予測の影響は小さくなり、現状の計算機の負荷だけでスケジュールするのと変わらない。

資源使用率が最も低い計算サーバに当該ジョブを割り当てる。

【0055】すなわち、この計算機 m は、ジョブ p 投入時に、資源 s について、投入前の実測値 LS_{ms} と、投入するジョブの資源要求量 LR_{ps} と、資源容量 LC_{ms} と、資源使用限界 B_{ms} とにおいて、

【0056】

【数6】

3)と合計することによってジョブ投入後の資源必要量を推定する。

【0060】すなわち、ステップ44で選択される計算機 m は、ジョブ p 投入時に、資源 s について、投入前の実測値 LS_{ms} と、投入するジョブの資源要求量 LR_{ps} と、資源容量 LC_{ms} と、資源使用限界 B_{ms} とにおいて、

【0061】

【数7】

時に、資源 s について、投入前の実測値 LS_{ms} と、投入するジョブの資源要求量 LR_{ps} と、資源容量 LC_{ms} と、資源使用限界 B_{ms} とにおいて、

【0065】

【数8】

【0072】繰り返し実行されるジョブが必要とするシステム資源量を自動的に獲得する。ジョブが必要とする資源量がわかると、投入後の計算機の資源使用が集中しないようにスケジュールできる。例えば、CPU負荷をかけるがIO負荷をかけないジョブと、CPU負荷をかけないがIO負荷をかけるジョブと、を組み合わせるスケジュールするようになり、計算機上のジョブ間の資源の奪い合いが少なくなり、効率的に計算機を使用できる。

【0073】

【発明の効果】個別のジョブが必要とする複数のシステム資源の量を自動的に獲得できる。この情報によって、特定のシステム資源に処理が集中することによって処理効率が落ちることを防ぐことができる。

【0074】個別ジョブの処理時間を推定できる。

【図面の簡単な説明】

【図1】本発明の実施例に係わる分散計算機システムの構成図である。

【図2】テーブルのテーブルの構成図である。

【図3】個別ジョブの資源要求量獲得方法のブロック図である。

【図4】ジョブ起動計算機割当スケジュール方法のブロック図である。

【符号の説明】

LRps …ジョブpが単位時間に必要な資源sの量の推定値、

LRms …計算機m上のジョブが単位時間に必要な資源sの量の推定値、

LApS …ジョブpが単位時間に使用する資源sの量の推定値、

LAmS …計算機mが単位時間に使用する資源sの量の推定値、

LSms …計算機mが単位時間に使用した資源sの量の実測値、

LCm …計算機mの単位時間に使用可能な資源sの容量、

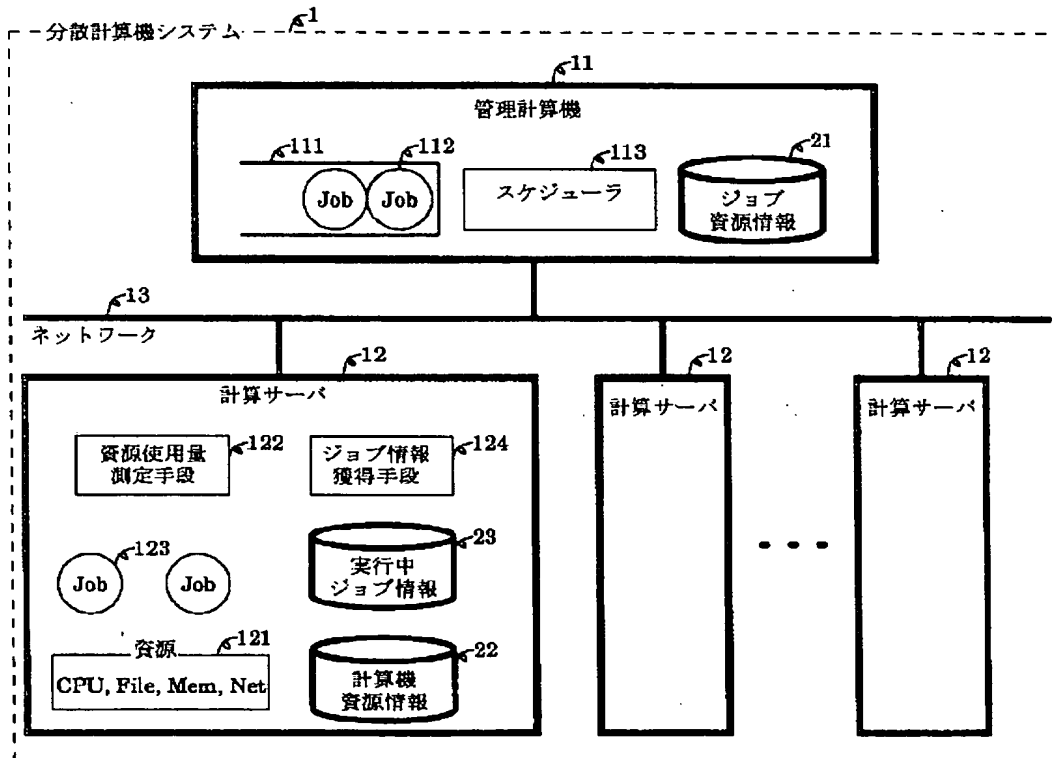
Eps …LApSの誤差、 ERps …LApSの誤差率、 $=Eps/LApS[\%]$ 、

Rps …LApSの確度、 $= (LApS - Eps) / LApS$ 、

Amp …計算機m上のジョブpの稼働の割合、Bms …資源使用限界。

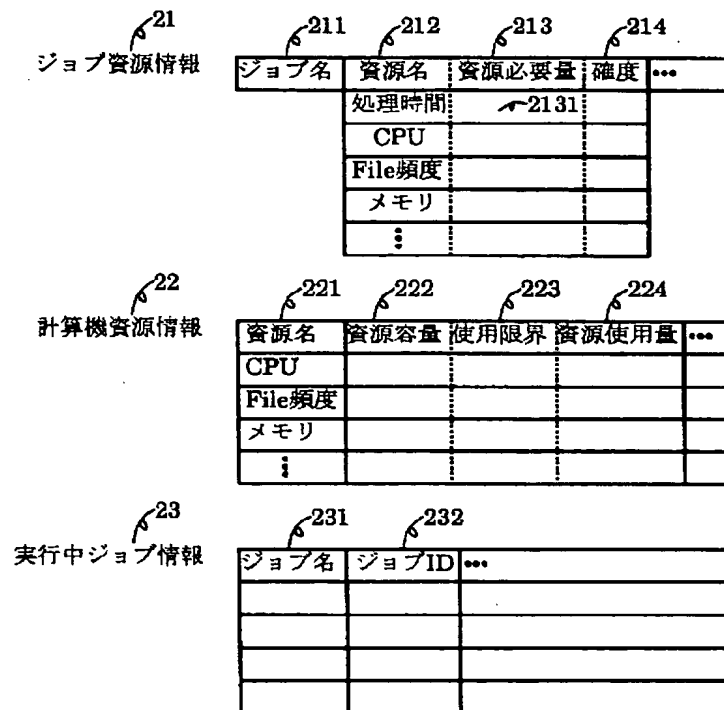
【図1】

図1



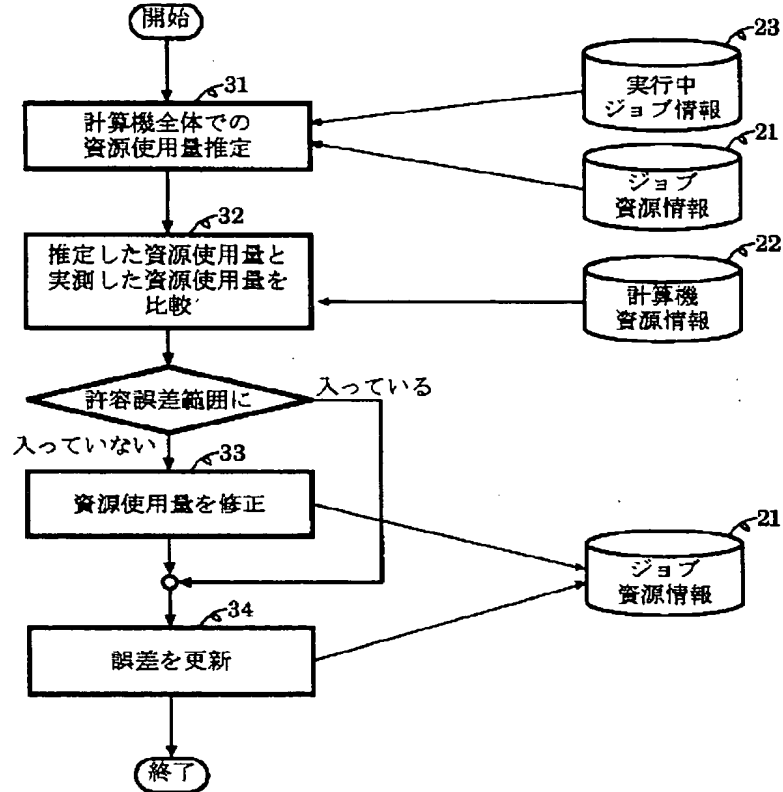
【図2】

図2



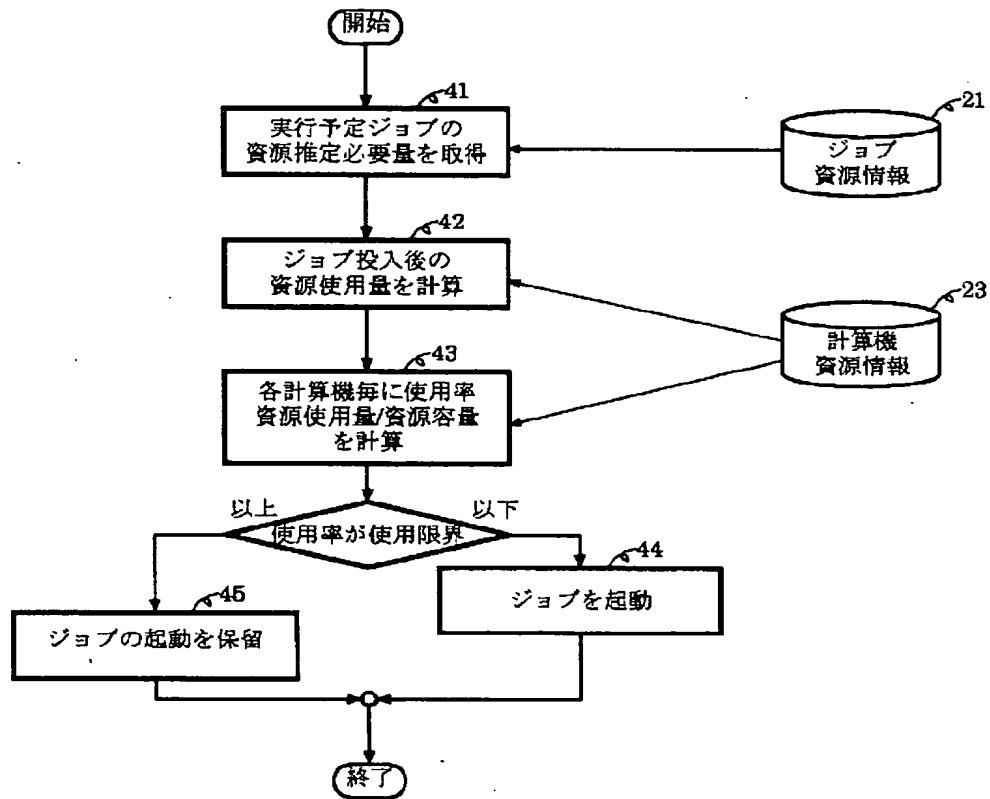
【図 3】

図 3



【図4】

図4



フロントページの続き

(72) 発明者 森田 真司
神奈川県横浜市戸塚区戸塚町5030番地株式
会社日立製作所ソフトウェア開発本部内